

Some New Performance Criteria in  
Robust Statistics-Small Sample Robustness

by

P. Papantoni-Kazakos

January 1977

Technical Report # 7701

SOME NEW PERFORMANCE CRITERIA IN  
ROBUST STATISTICS -- SMALL SAMPLE ROBUSTNESS\*

by

P. Papantoni-Kazakos  
Department of Electrical Engineering  
Rice University

ABSTRACT

In this paper the analysis of estimates operating on dependent data is considered. Special dependent data structures are considered and the analysis is made for three different choices of contamination and performance distance measures.

For contamination and performance measures both being Lévy (Hampel model), an analysis that is particularly oriented toward fast convergence of the estimate to a value that is stable (robust) inside the contaminated family is undertaken. The minimum sample size to satisfy certain performance is investigated and it is found that the problem reduces to finding continuous, absolutely bounded estimates with logarithms of their moment generating function slowly increasing with the absolute value of the argument for all data distributions considered.

1. INTRODUCTION

Humans have always taken satisfaction in the game of "outsmarting Nature." This game becomes particularly intriguing when Nature tries to confuse its human opponents by slightly changing the underlying statistical rules that specify its behavior. To beat Nature's game, then, the human player must follow strategies that are not very sensitive to such statistical changes.

---

\* This work was supported by the Air Force Grant AFOSR 77-3156.

The design of such statistical variation-resistant strategies varies with the particular objective of the human player at the time and the decision rules evolving from this design are called robust.

Robust decision rules have been studied by several statisticians [1,3,7,10,12] and engineers [11,13,18] from different points of view. The promising qualitative analysis of robust estimation first introduced by Hampel [10] has been used and extended [12,18], but the performance criterion has always been expressed through the Lèvy distance of the data as well as the estimator distributions. Specifically, an estimate  $\hat{s}(X_n)$  using the data  $X_n$  has been called robust (or weakly robust according to [18]) at some distribution  $Q_0$ , if every distribution  $Q$  that is close to  $Q_0$  in the Lèvy distance sense results to  $\hat{s}(X_n)$  distribution that is Lèvy distance close to the  $\hat{s}(X_n)$  distribution resulting from  $Q_0$ -distributed data. This definition of robustness, while leading to constructive analysis of robust estimators [10,18], has two disadvantages: It does not offer convergence rates of the estimates for large number of data and it may be too demanding or even not representative enough for some applications. In addition, the constructive analysis of "Lèvy robust" estimators has been accomplished only for independent  $X_n$  components.

In particular, one of the properties that characterize an estimate  $\hat{s}(X_n)$  that is "Lèvy robust" at some one-dimensional distribution  $F_0$  and is applied on independent data, is continuity at  $F_0$  [10,18]. This property, which actually means closeness of the values  $\hat{s}(X_n)$ ,  $\hat{s}(X_m)$  for vectors  $X_n$ ,  $X_m$  that specify experimental distributions that are both close to  $F_0$ , guarantees convergence of  $\hat{s}(X_n)$  to some constant depending on  $F_0$ , but it does not specify the convergence rate.

Also, in some applications the preservation of some specific, less general than the Lèvy distance, characteristic may be desirable. In this case the per-

formance criterion is different and so are the desirable properties of the estimator. Performance criteria that are easier to calculate than the Lévy distance may be preferred.

Finally, the design of estimators that are "robust" in the presence of dependent data structures is certainly a problem that is challenging as well as realistically interesting.

In section 2 of the present paper some preliminary discussion on previous analysis of robustness on statistically contaminated distribution and on possible spaces, different performance criteria is presented.

In section 3, the design of "Vasarshtein robust" estimators is undertaken. The dependence structure of the data is naturally incorporated in this case if the distortion or penalty measure is square difference.

In section 4, properties of estimators that are robust as mappings from a data Lévy-contaminated space to an estimate space characterized by a Vasarshtein performance criterion are discussed.

In section 5, a design method of "Lévy robust estimates" is presented that incorporates an exponential convergence rate. Special dependence structure of the data is considered.

Section 6 includes examples of estimators that are "robust" in the senses of sections 3, 4 and 5.

## 2. PRELIMINARIES

Robustness has lately been defined [10,18] as stability of some stochastic distance measure defined on the estimator probability space. Specifically, if  $X_n$  denotes the vector of  $n$  discrete data,  $Q_0$  is some well-known multi-dimensional cumulative distribution applied on  $X_n$  (where  $Q_{0n}$  is the  $n$ -dimensional distribution evolving from  $Q_0$ ),  $Q$  is some arbitrary cumulative distribution on  $X_n$ ;  $n = 1, 2, \dots$ ;  $d_1(\cdot, \cdot)$  is a stochastic distance measure

defined on the data distribution space;  $\hat{s}_n(X_n)$  is a scalar estimate that is a deterministic function of the data  $X_n$ ;  $D(\hat{s}_n)$ ,  $D^0(\hat{s}_n)$  are the distributions of  $\hat{s}_n(X_n)$  determined through  $Q$ ,  $Q^0$  respectively; and  $d_{2n}(\cdot, \cdot)$  is a stochastic distance measure defined on the distribution space  $D(\hat{s}_n)$ , then the sequence  $\{\hat{s}_n\}$  is weakly robust at  $Q^0$  (as defined in [18]) if given  $\epsilon > 0$ , there is some  $\delta(\epsilon) > 0$  such that: For every  $Q$  satisfying  $d_1(Q^0, Q) < \delta(\epsilon)$ ,  $d_{2n}(D^0(\hat{s}_n), D(\hat{s}_n)) < \epsilon$ ;  $\forall n$  is implied.

It is obvious from the above definition of weak robustness that the stochastic distances  $d_1(\cdot, \cdot)$ ,  $d_{2n}(\cdot, \cdot)$  must be chosen carefully to satisfy the designer's specific objective. In particular, the distance  $d_1(\cdot, \cdot)$  represents the kind of  $Q^0$  statistical contamination considered, while  $d_{2n}(\cdot, \cdot)$  is the performance measure of the estimate.

In the choice of  $d_1(\cdot, \cdot)$  and  $d_{2n}(\cdot, \cdot)$  the good representation of the particular model as well as the calculability of the distances must be taken into consideration.

In the study presented in [17], it became apparent that among the plethora of stochastic distances existing, there are some more and some less approachable. Specifically, the Lèvy distance although useful in the robust analysis is often very hard to calculate. On the other hand, the Vasershtein distance is simpler in some cases and if also representative of the model considered, it becomes an excellent contamination or performance measure choice.

In the present study, only Lèvy and Vasershtein distance choices will be considered. Therefore, their definition and some of their properties that are related to the analysis in this paper are presented.

Both Lèvy and Vasershtein distances include a distortion or penalty measure  $\rho(\cdot, \cdot)$  applied on the outcomes of the distributions involved. Specifically, if we concentrate our attention on discrete data structures, let  $X_n$ ,

$Y_n$  be two different  $n$ -data values and let  $\rho(X_n, Y_n)$  be their relative distortion. Also, let  $Q_n^0, Q_n$  be two different cumulative distributions of  $X_n$ . Then, the Lèvy distance  $d_{L\rho}(Q_n^0, Q_n)$  and the Vasershtein distance  $d_{V\rho}(Q_n^0, Q_n)$  are defined as follows:

$$d_{L\rho}(Q_n^0, Q_n) = \inf\{\epsilon: Q_n^0(X_n) \leq Q_n(UY_n: \rho(X_n, Y_n) \leq \epsilon) + \epsilon, \\ Q_n(X_n) \leq Q_n^0(UY_n: \rho(X_n, Y_n) \leq \epsilon) + \epsilon; \forall X_n\} \quad (1)$$

$$d_{V\rho}(Q_n^0, Q_n) = \inf E_{Q_{2n}(\cdot, \cdot)}\{\rho(X_n, Y_n)\} \quad (2)$$

all  $Q_{2n}(\cdot, \cdot)$  inducing  
 $Q_n^0, Q_n$  marginals

The same distances can be defined for general multivariate distributions  $Q^0, Q$ , inducing  $Q_n^0, Q_n \forall n$  in the following way:

$$d_{L\rho}(Q^0, Q) = \sup_n d_{L\rho}(Q_n^0, Q_n) \quad (3)$$

$$d_{V\rho}(Q^0, Q) = \sup_n d_{V\rho}(Q_n^0, Q_n) \quad (4)$$

In [18] it was shown that  $d_{L\rho}(Q_n^0, Q_n)$  is nondecreasing with increasing  $n$  for  $Q_n^0, Q_n$  being both products of one-dimensional distributions. It is straight forward to show that the same is true for arbitrary  $Q_n^0, Q_n$ . The Vasershtein distance  $d_{V\rho}(Q_n^0, Q_n)$ , on the other hand, can be independent of  $n$  for stationary data structures and proper choice of the distortion (penalty) measure  $\rho(\cdot, \cdot)$ . Indeed, let  $X_n = \{x_i; i = 1, \dots, n\}$  be a sample vector from a wide-stationary process  $x(t)$  whose autocovariance function is  $R^0(\tau)$  and whose  $n$ -dimensional discrete distribution is represented by  $Q_n^0$ . Also, let  $Q_n$  be the discrete representation of another wide stationary process  $y(t)$  with autocovariance function  $R(\tau)$ . If the Vasershtein distance is defined in the space of jointly stationary distributions and if

$\rho(X_n, Y_n) = \frac{1}{n} [X_n - Y_n]' [X_n - Y_n]$ , where  $[ ]'$  means transpose, then as it was found in [15] and [17], one obtains

$$d_{V_p}(Q_n^0, Q_n) = \inf_{R^C(0)} \{R^0(0) + R(0) - 2R^C(0) + (m^0 - m)^2\} \quad (5)$$

where  $R^C(\tau)$  is the crosscovariance function of  $x(t), y(t)$ ,  $m^0$  is the mean of  $x(t)$  and  $m$  the mean of  $y(t)$ . The expression in (5) is obviously independent of  $n$ . Furthermore, if

$$P^0(\lambda) = \sum_{k=-\infty}^{\infty} R^0(k) e^{jk\lambda} \Rightarrow R^0(0) = (2\pi)^{-1} \int_{-\pi}^{\pi} P^0(\lambda) d\lambda \quad (6)$$

$$P(\lambda) = \sum_{k=-\infty}^{\infty} R(k) e^{jk\lambda} \Rightarrow R(0) = (2\pi)^{-1} \int_{-\pi}^{\pi} P(\lambda) d\lambda \quad (7)$$

if  $Q_n, Q_n^0$  are Gaussian, then the  $R^C(0)$  that satisfies the infimum in (5) was found in [15] to be given by the expression:

$$R^C(0) = (2\pi)^{-1} \int_{-\pi}^{\pi} \sqrt{P^0(\lambda) P(\lambda)} d\lambda \quad (8)$$

which provides the following expression of the Vasershtein distance for stationary distribution Gaussian spaces:

$$d_{V_p}(Q_n^0, Q_n) = d_{V_p}(Q^0, Q) = (2\pi)^{-1} \int_{-\pi}^{\pi} (\sqrt{P^0(\lambda)} - \sqrt{P(\lambda)})^2 d\lambda + (m^0 - m)^2 \quad (9)$$

$$\text{where } \rho(X_n, Y_n) = \frac{1}{n} [X_n - Y_n]' [X_n - Y_n] \quad (10)$$

If  $Q_n^0, Q_n$  are not Gaussian, the expression on the left of (9) becomes a lower bound on  $d_{V_p}(Q_n^0, Q_n)$ .

The expression in (9) is simple, computable analytically in most cases, and independent of the sample size  $n$ . It is also significant to observe from (9) that the wide sense stationary Gaussian distributions that are "close" with respect to a square error, are the ones whose means and discrete spectra functions are "close."

The constructive analysis of "robust" estimators that was presented in [10] and [18] in addition to being "Lèvy robust," it was also limited in the case of independent data structures. It was found there that the sequence  $\{\hat{s}_n(X_n)\}$  of estimates is weakly robust at some  $Q_{01}$  in the sense of both the  $d_1(\cdot, \cdot)$ ,  $d_{2n}(\cdot, \cdot)$  distances being Lèvy, if: a)  $\hat{s}_n(X_n)$  is continuous for every  $n$  as a real function with  $E^n$  Euclidean domain, b)  $\hat{s}_n(X_n)$  is continuous at  $Q_{01}$ , this continuity meaning that for every  $X_n, X_m$  such that they determine experimental distributions  $n_{X_n}(y), n_{Y_m}(y)$  that are both Lèvy close to  $Q_{01}$ , it is implied that  $|\hat{s}_n(X_n) - \hat{s}_m(Y_m)|$  is small.

The above analysis does not provide convergence rates of the estimate  $\hat{s}_n(X_n)$ . Such a rate is important to the designer that is looking not only for "robust" estimates but also for sufficient sample sizes also. It is important to have an analysis that answers the double question: "What kind of estimate will be robust for a given contaminated family and how many data are sufficient to guarantee a certain minimum level of performance inside the same family?"

An attempt to answer the above question for certain dependent data structures is even more valuable. Having the dependent data situations in mind, we will first study estimates that are robust for contamination and performance measures that are not both Lèvy. Then, we will present an analysis that although applied to Lèvy distribution contamination and Lèvy performance measure, it defers from the classical one in [10] and [18] in the fact that it can apply to dependent data and it incorporates the convergence rate of the estimates.

### 3. VASERSHTEIN ROBUST ESTIMATORS

In this section we will assume that the human player has the information that Nature uses a Vasershtein algorithm to contaminate its underline statistics.



That is, we consider the measure of contamination to be the Vasershtein distance and we pick as distortion measure the square error one. If, in addition, Nature picks its statistics from the wide-sense stationary distribution space  $\mathcal{F}$  that surrounds a well-known distribution  $Q^0$ , then the Vasershtein distance between  $Q_0$  and an arbitrary member  $Q \in \mathcal{F}$  is given by (as expressed in (9) of section 2):

$$d_{V_p}(Q^0, Q) \geq (2\pi)^{-1} \int_{-\pi}^{\pi} (\sqrt{P^0(\lambda)} - \sqrt{P(\lambda)})^2 d\lambda + (m^0 - m)^2 \quad (11)$$

where

$$\rho(X_n, Y_n) = \frac{1}{n} [X_n - Y_n]' [X_n - Y_n] \quad (12)$$

$m^0, m$  the one-dimensional means corresponding to  $Q^0$  and  $Q$  and  $P^0(\lambda)$ ,  $P(\lambda)$  the respective discrete spectral densities. We have equality in (11) if  $Q_n^0, Q_n$  distributions are both Gaussian.

From (11) it is apparent that the Vasershtein  $\rho$ -contaminated Gaussian distribution families are the families with contaminated means and spectral densities.

For consideration of arbitrary  $Q$  distributions (even when  $Q^0$  is Gaussian) we can transform the contamination measure to:

$$d_{P,m}(Q^0, Q) = (2\pi)^{-1} \int_{-\pi}^{\pi} (\sqrt{P^0(\lambda)} - \sqrt{P(\lambda)})^2 d\lambda + (m^0 - m)^2 \quad (13)$$

In other words, we suppose that nature contaminates the data statistics by contaminating the spectral density and the mean of the underline stationary process.

Let us now suppose that the observer's performance measure is the mean square one. That is, if  $p$  parameters must be estimated from the collected data  $X_n$ , then the robustness of the  $p$ -dimensional estimate  $\hat{S}_n(X_n)$  is evaluated through a mean square error value. In other words, the estimate designer is fully satisfied if the average mean square distortion between

the  $\hat{S}_n(X_n)$  he calculates when Nature uses  $Q^0$  underline statistics and the value  $\hat{S}_n(Y_n)$  he finds when  $Q$  is true, remains small when  $Q^0, Q$  are close enough in the Vasershtein distance sense.

If  $D^0(\hat{S}_n), D(\hat{S}_n)$  are the distributions of  $\hat{S}_n(X_n)$  evolving from  $Q^0, Q$  respectively, then

$$d_{V_p}(D^0(\hat{S}_n), D(\hat{S}_n)) = p^{-1} \sum_{i=1}^p \inf_{r_{ni}(\cdot, \cdot) \text{ inducing } D^0(\hat{s}_{in}), D(\hat{s}_{in})} E_{r_{ni}(\cdot, \cdot)} \{ \hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n) \}^2 \quad (14)$$

where

$$\hat{S}_n(X_n) = \{ \hat{s}_{in}(X_n); i = 1, \dots, p \} \quad (15)$$

If

$$E_{D^0(\hat{s}_{in})} \{ \hat{s}_{in}(X_n) \} = m_{in}^0; E_{D(\hat{s}_{in})} \{ \hat{s}_{in}(X_n) \} = m_{in} \quad (16)$$

$$\begin{aligned} E_{D^0(\hat{s}_{in})} \{ [\hat{s}_{in}(X_n) - m_{in}^0]^2 \} &= \sigma_{in}^0(0); E_{D(\hat{s}_{in})} \{ [\hat{s}_{in}(X_n) - m_{in}]^2 \} \\ &= \sigma_{in}(0) \end{aligned} \quad (17)$$

The following lemma can be expressed:

#### Lemma 1

The distance in (14) is bounded from below by the expression

$$d_{\sigma, m}(D^0(\hat{S}_n), D(\hat{S}_n)) = \sum_{i=1}^p \{ (\sqrt{\sigma_{in}^0(0)} - \sqrt{\sigma_{in}(0)})^2 + (m_{in}^0 - m_{in})^2 \} \quad (18)$$

#### Proof

For any joint distribution  $r_{ni}(\cdot, \cdot)$  with  $D^0(\hat{s}_{in}), D(\hat{s}_{in})$  marginals, if  $\sigma_{in}^i(0)$  the crosscovariance determined by  $r_{ni}(\cdot, \cdot)$ , then, the matrix

$$\begin{bmatrix} \sigma_{in}^o(0) & \sigma_{in}^c(0) \\ \sigma_{in}^c(0) & \sigma_{in}(0) \end{bmatrix}$$

must be nonnegative definite.

Therefore,

$$\sqrt{\sigma_{in}^o(0) \sigma_{in}(0)} \geq \sigma_{in}^c(0)$$

Then,

$$\begin{aligned} E_{r_{ni}}(\cdot, \cdot) \{ \hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n) \}^2 &= \sigma_{in}^o(0) + \sigma_{in}(0) - 2\sigma_{in}^c(0) \\ &+ (m_{in}^o - m_{in})^2 \geq \sigma_{in}^o(0) + \sigma_{in}(0) - 2\sqrt{\sigma_{in}^o(0) \sigma_{in}(0)} + (m_{in}^o - m_{in})^2 \\ &= [\sqrt{\sigma_{in}^o(0)} - \sqrt{\sigma_{in}(0)}]^2 + (m_{in}^o - m_{in})^2 \end{aligned}$$

and

$$d_V^o(D^o(\hat{S}_n), D(\hat{S}_n)) \geq p^{-1} \sum_{i=1}^p \{ [\sqrt{\sigma_{in}^o(0)} - \sqrt{\sigma_{in}(0)}]^2 + (m_{in}^o - m_{in})^2 \}$$

The expression in (18) is equal to the distance in (14) if and only if  $D(\hat{s}_{in})$  can be the distribution of a linear transformation of the variable distributed as in  $D^o(\hat{s}_{in})$ , for every  $1 \leq i \leq p$ .

We will summarize the observations we made up to now by the following three definition:

#### Definition 1

We will call a sequence  $\{\hat{S}_n\}$  of  $p$ -dimensional estimates  $p$ -Vasershtein weakly robust inside a wide sense stationary distribution family  $\mathcal{F}$  and at some  $Q^o$  if and only if given  $\epsilon > 0$ , there is some  $\delta(\epsilon) > 0$  such that: For  $P^o(\lambda), m^o$  being the spectral density and the mean induced by  $Q^o$ ,  $P(\lambda), m$  being the spectral density and mean of some  $Q \in \mathcal{F}$  and for

$$(2\pi)^{-1} \int_{-\pi}^{\pi} (\sqrt{P^0(\lambda)} - \sqrt{P(\lambda)})^2 d\lambda + (m^0 - m)^2 < \delta(\epsilon)$$

it is implied that:

$$p^{-1} \sum_{i=1}^p \{[\sqrt{\sigma_{in}^0(0)} - \sqrt{\sigma_{in}(0)}]^2 + (m_{in}^0 - m_{in})^2\} < \epsilon; \forall n$$

where  $\sigma_{in}^0(0)$ ,  $\sigma_{in}(0)$ ,  $m_{in}^0$ ,  $m_{in}$  are given by (17) and (16) respectively.

### Observations

1. If the  $\mathcal{F}$  family in definition 1 is a Gaussian stationary family and the estimates  $\hat{s}_{in}(X_n)$  are linear transformations of the data, then the expression

$$(2\pi)^{-1} \int_{-\pi}^{\pi} (\sqrt{P^0(\lambda)} - \sqrt{P(\lambda)})^2 d\lambda + (m^0 - m)^2$$

$$p^{-1} \sum_{i=1}^p \{(\sqrt{\sigma_{in}^0(0)} - \sqrt{\sigma_{in}(0)})^2 + (m_{in}^0 - m_{in})^2\}$$

are the exact Vasershtein distances of the distributions  $Q^0$ ,  $Q$  and  $D^0(\hat{s}_{in})$ ,  $D(\hat{s}_{in})$ , respectively.

2. It is evident from definition 1 that since we want the closeness of the estimate means and variances guaranteed by the closeness of just the spectral densities and the means of the data distributions, we must limit the estimates to linear transformations of the observations.

Let us define the linear estimates

$$\hat{s}_{in}(X_n) = \sum_{k=1}^n a_{ni}(k) x_k \quad (19)$$

where  $X_n = \{x_k; k = 1, \dots, n\}$  and the coefficients  $a_{ni}(k)$  are real, scalar, and, in general, different for different dimensionality  $n$ . To study robustness, we need the means and variances of the estimate in (19) under data

distributions  $Q^0$  and  $Q$ . Indeed, we have

$$\begin{aligned}
 m_{in}^0 &= E\{\hat{s}_{in}(X_n)/Q^0 \text{ distr}\} = m^0 \sum_{k=1}^n a_{ni}(k) \\
 m_{in} &= E\{\hat{s}_{in}(X_n)/Q \text{ distr}\} = m \sum_{k=1}^n a_{ni}(k) \\
 \sigma_{in}^0(0) &= \sum_{k,\ell=1}^n a_{ni}(k) a_{ni}(\ell) E\{(x_k - m^0)(x_\ell - m^0)\} \\
 &= \sum_{k,\ell=1}^n a_{ni}(k) a_{ni}(\ell) R^0(k - \ell) \\
 &= (2\pi)^{-1} \int_{-\pi}^{\pi} P^0(\lambda) \sum_{k,\ell=1}^n a_{ni}(k) a_{ni}(\ell) e^{-j(k-\ell)\lambda} d\lambda
 \end{aligned} \tag{20}$$

where  $Q^0$  is wide sense stationary with autocovariance  $R^0(\tau)$  and power spectral density  $P^0(\lambda)$ .

From the above expression we finally obtain:

$$\sigma_{in}^0(0) = (2\pi)^{-1} \int_{-\pi}^{\pi} P^0(\lambda) \left\| \sum_{k=1}^n a_{ni}(k) e^{-jk\lambda} \right\|^2 d\lambda \tag{21}$$

and similarly:

$$\sigma_{in}(0) = (2\pi)^{-1} \int_{-\pi}^{\pi} P(\lambda) \left\| \sum_{k=1}^n a_{ni}(k) e^{-jk\lambda} \right\|^2 d\lambda \tag{22}$$

for the variance under distribution  $Q$ . Applying the Schwartz inequality:

$$\begin{aligned}
 \int_{-\pi}^{\pi} \sqrt{P^0(\lambda)} \sqrt{P(\lambda)} \left\| \sum_{k=1}^n a_{ni}(k) e^{-jk\lambda} \right\|^2 d\lambda &\leq \left[ \int_{-\pi}^{\pi} P^0(\lambda) \left\| \sum_{k=1}^n a_{ni}(k) e^{-jk\lambda} \right\|^2 d\lambda \right]^{\frac{1}{2}} \cdot \\
 &\quad \cdot \left[ \int_{-\pi}^{\pi} P(\lambda) \left\| \sum_{k=1}^n a_{ni}(k) e^{-jk\lambda} \right\|^2 d\lambda \right]^{\frac{1}{2}}
 \end{aligned}$$

on  $[\sqrt{\sigma_{in}^0(0)} - \sqrt{\sigma_{in}(0)}]^2$

we obtain from (21) and (22):

$$[\sqrt{\sigma_{in}^o(0)} - \sqrt{\sigma_{in}(0)}]^2 \leq (2\pi)^{-1} \int_{-\pi}^{\pi} [\sqrt{P^o(\lambda)} - \sqrt{P(\lambda)}]^2 \left\| \sum_{k=1}^n a_{ni}(k) e^{-jk\lambda} \right\|^2 d\lambda \quad (23)$$

while from (21) we get directly:

$$(m_{in}^o - m_{in})^2 = (m^o - m)^2 \left[ \sum_{k=1}^n a_{ni}(k) \right]^2 \quad (24)$$

The expressions in (23) and (24) express the connection between data and estimate statistics needed to specify  $\rho$ -Vasershtein weak robustness, as given by definition 1.

From the analysis done above, and the expression of  $\rho$ -Vasershtein weak robustness in definition 1, a lemma offering a constructive properties of estimates that are Vasershtein weak robust is obtained.

#### Lemma 2

A  $p$ -dimensional estimate  $\hat{S}_n(X_n)$  that is  $\rho$ -Vasershtein weak robust, as expressed by definition 1, must be linear. If this linear estimate is given by the expressions:

$$\hat{S}_n(X_n) = \{\hat{s}_{in}(X_n); i = 1, \dots, p\}, \quad \hat{s}_{in}(X_n) = \sum_{k=1}^n a_{ni}(k) x_k$$

a sufficient condition for the present sense of robustness is that the sequences

$$\left\{ \sum_{k=1}^n |a_{ni}(k)| \right\}$$

converge to some finite value for every  $1 \leq i \leq p$ .

#### Proof:

Let the sequence  $\left\{ \sum_{k=1}^n |a_{ni}(k)| \right\}$  converge to some  $A_i < \infty$ . Then, also

$$\left( \sum_{k=1}^n a_{ni}(k) \right)^2 \leq \left( \sum_{k=1}^n |a_{ni}(k)| \right)^2 \leq A_i^2 \leq \max_{1 \leq i \leq p} A_i^2; \quad \forall n.$$

Given  $\epsilon > 0$ , pick

$$\delta(\epsilon) = \frac{\epsilon}{2p \cdot \max_{1 \leq i \leq p} A_i^2}$$

and the conditions in definition 1 are satisfied.

---

According to lemma 2, the experimental mean  $\hat{s}_{in}(X_n) = \frac{1}{n} \sum_{i=1}^n x_i$  is a  $p$ -Vasershtein weakly robust estimate. The condition of lemma 2 can be easily seen to satisfy robustness (not just weak robustness) properties inside some data-distribution contaminated family.

Concluding this section, we want to emphasize that the robustness structure presented here considers dependent data with dependence expressed by arbitrary wide-sense stationary distributions. This dependence was explicitly incorporated in the robustness only through the spectral densities of these distributions.

#### 4. LÉVY-VASERSHTEIN ROBUST ESTIMATORS

In this section we consider the case that the contamination measure on the data distribution space is the Lévy distance, while the performance measure on the space of the estimates is the Vasershtein-type distance

$$d_{vt}(D^0(\hat{S}_n), D(\hat{S}_n)) = p^{-1} \sum_{i=1}^p \{ [\sqrt{\sigma_{in}^0(0)} - \sqrt{\sigma_{in}(0)}]^2 + (m_{in}^0 - m_{in})^2 \} \quad (25)$$

The characteristics  $\sigma_{in}^0(0)$ ,  $\sigma_{in}(0)$ ,  $m_{in}^0$ ,  $m_{in}$  are given by expressions (16) and (17) of the previous section and the distance in (25) is equal to the Vasershtein distance again if the data distribution family is a wide sense stationary family and the estimates are linear.

The robustness considered here is precisely expressed by the following definition.

Definition 2

A sequence of  $p$ -dimensional estimates  $\{\hat{S}_n(X_n)\}$  is weakly  $\rho$ -Lévy-Vasershtein robust at some distribution  $Q^0$  if and only if: given  $\epsilon > 0$ , there is some  $\delta(\epsilon) > 0$  such that for every distribution  $Q$  satisfying:

$$d_{L_p}(Q^0, Q) < \delta(\epsilon)$$

it is implied that

$$d_{V_t}(D^0(\hat{S}_n), D(\hat{S}_n)) < \epsilon ; \forall n .$$

---

The distance  $d_{L_p}(Q^0, Q)$  is defined by expressions (1) and (3) and the distance  $d_{V_t}(D^0(\hat{S}_n), D(\hat{S}_n))$  by expression (25). The dependence structure of the data (as expressed by  $Q^0$  and  $Q$ ) is arbitrary at this point.

In the analysis for the discovery of constructive properties of the estimates that are robust in the  $\rho$ -Lévy-Vasershtein sense, we will need to bound the absolute values of the  $\hat{S}_n(X_n)$  components  $\hat{s}_{in}(X_n)$ ;  $i = 1, \dots, p$ , for every  $i$  and every  $n$ . That restriction is mostly very useful realistically whenever we are seeking the estimation of parameters whose values (we know in advance) move inside a limited interval. The value restriction on the estimates rejects a priori the unacceptably (or dangerously) false decisions on the parameters of interest.

Similarly to the method presented in [18], we will break the constructive analysis of the  $\rho$ -Lévy-Vasershtein weakly robust estimates into two parts. One for sample sizes  $n$  bounded from above by some  $n_0$  and one for the  $n$ 's that exceed this bound  $n_0$ .

We proceed first with the bounded  $n$  part, presenting the following lemma.

Lemma 3

Let an estimate  $\hat{S}_n(X_n) = \{\hat{s}_{in}(X_n); 1 \leq i \leq p\}$  be absolutely bounded for



every  $i, n, X_n$  and let all the  $\hat{s}_{in}(X_n)$  components be continuous as real functions defined on the  $E^n$  Euclidean space  $V_n$ . Then, given some natural number  $n_0$  and some  $\epsilon > 0$ , there is some  $\delta(\epsilon, n_0) > 0$  such that  $\forall n \leq n_0$  and for  $Q_n$ :

$$d_{L_p}(Q_n^0, Q_n) < \delta(\epsilon)$$

it is implied that:

$$d_{V_t}(D^0(\hat{S}_n), D(\hat{S}_n)) < \epsilon$$

The proof of this lemma is presented in appendix A. The continuity of  $\hat{s}_{in}(X_n)$  as a real function is from the  $\rho(X_n, Y_n)$  distortion measure on the data that is incorporated in the Lèvy distance, to the absolute value difference  $|\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)|$  of the estimates. No consideration of particular dependence structure of the data was necessary at this point.

The lemma we present next combines the properties of the estimator that satisfy the  $\rho$ -Lèvy-Vasershtein weak robust requirements for finite as well as infinite sample sizes  $n$ .

For the transition to the infinite  $n$  step, the specification of a particular dependence structure of the data is necessary. In particular, we will assume that the family of data statistics considered is limited to  $m$ -dependent distributions. In addition, the data will be collected in groups of  $k$  consecutive data and the groups will be in distance of  $m$  data from each other. Specifically, the data vector  $X_n$  will consist, in this case, of  $k$ -dimensional vectors  $X_{ki}; i = 1, 2, \dots, n_1$ . The components of each  $X_{ki}$  vector are dependent, but  $X_{ki}$  is independent of  $X_{kj}$  for  $i \neq j$ .

The experimental distribution of the vector  $X_n$  is defined then as follows

$$n_{X_n}^k(\omega_1, \dots, \omega_k) = \frac{1}{n_1} \{ \# \text{ of } X_{ki} \text{'s with } x_1 < \omega_1, x_2 < \omega_2, \dots, x_k < \omega_k \} \quad (26)$$

A continuity of the estimate that is actually a property of the estimator indicating stochastic stability around the data central distribution  $Q^0$  is defined below:

### Definition 3

For data  $X_n$  consisting of  $X_{ki}$  independent vectors and experimental distributions defined by (26), an estimator  $\hat{S}_n(X_n) = \{\hat{s}_{in}(X_n); 1 \leq i \leq p\}$  is continuous at  $Q_k^0$  if and only if:

Given  $\epsilon > 0$ , there is some  $\mu(\epsilon) > 0$ ,  $n_0$  such that: For every  $X_n$ ,  $Y_n$  satisfying

$$d_{L_\rho}(n_{X_n}^k, Q_k^0) < \mu(\epsilon)$$

$$d_{L_\rho}(n_{Y_n}^k, Q_k^0) < \mu(\epsilon)$$

and  $n > n_0$

it is implied:

$$\max_i |\hat{s}_{in}(X_n) - \hat{s}_{im}(Y_n)| < \epsilon$$

Combination of lemma 3 and definition 3 leads to the following lemma whose proof can be found in Appendix A.

### Lemma 4

Let  $\hat{S}_n(X_n) = \{\hat{s}_{in}(X_n); 1 \leq i \leq p\}$  be absolutely bounded for every  $i$ ,  $n$ ,  $X_n$ . Let every  $\hat{s}_{in}(X_n)$  be a  $\rho$ -continuous real function on  $E^n \forall n$  and let  $\hat{S}_n(X_n)$  be continuous at  $Q_k^0$  where  $X_n$  is formed from  $n_1$  independent  $k$ -dimensional data vectors. Then the estimate  $\hat{S}_n(X_n)$  is  $\rho$ -Lèvy-Vasershtein weakly robust.

The  $\rho$  included in lemma 4 and definition 2 indicates the distortion measure included in the Lèvy distance while, as we show in the previous section, the Vasershtein type performance  $d_{V_t}(\cdot, \cdot)$  is the result of the consideration of the square error data measure in the expression of the Vasershtein distance.

The estimator properties that guarantee weak robustness through the contamination and performance measures considered in this section are similar to the ones of the classical Lèvy-Lèvy model analyzed in [10] and [18]. Here the absolute boundness of the estimate is an additional desirable property.

The estimators that are not weakly robust in the Lèvy-Lèvy sense are not robust in the Lèvy-Vasershtein sense also.

The means for the design of the properly "robust estimates are similar in both of the above cases.

In the following section we present an alternative analysis method that incorporates convergence rates and gives us a better feeling as to the proper design methods for Lèvy contamination, Lèvy performance robust estimates.

## 5. A NEW APPROACH TO LÈVY-LÈVY ROBUST ESTIMATION

The robust models that were considered by Hampel [10] and Papantoni-Kazakos [18] were based on Lèvy-contaminated data distribution families and Lèvy performance criterion of the estimates. Specifically, if  $\rho(X_n, Y_n)$  is some distortion (penalty) measure defined on the data, robustness is defined as follows according to this model:

### Definition 4

A sequence  $\{\hat{S}_n(X_n)\}$  of  $p$ -dimensional estimates is Lèvy-Lèvy weakly robust at  $Q^0$  if and only if given  $\epsilon > 0$ , there is some  $\delta(\epsilon) > 0$  such that:

For every  $Q$  such that

$$d_{L_p}(Q^0, Q) < \delta(\epsilon)$$

it is implied that:

$$d_{L_p}(D^0(\hat{S}_n), D(\hat{S}_n)) < \epsilon ; \forall n$$

---

$D^0(\hat{S}_n)$ ,  $D(\hat{S}_n)$  are the estimate  $p$ -dimensional distributions induced by  $Q^0(X_n)$ ,  $Q(X_n)$ , respectively. Also, the distributions  $Q^0$ ,  $Q$  generate, in general, dependent  $X_n$  vectors.

If  $n_0$  is some finite natural number, the stability property expressed by definition 4 is satisfied for  $n \leq n_0$  if the components  $\hat{s}_{in}(X_n)$  of the estimate  $\hat{S}_n(X_n)$  are all continuous as functions with  $E^n$  Euclidean domain, where  $n$  any natural number. The proof of this is similar to the corresponding proof for independent data appearing in [18] and to the proof of lemma 3 of the present paper that can be found in appendix A. Formally speaking:

#### Lemma 4

If  $\hat{S}_n(X_n)$  is continuous as a function on  $E^n \forall n$ , then, given  $\epsilon > 0$  and  $n_0$ , there is some  $\delta(\epsilon, n_0) > 0$  such that:

For every distribution  $Q$  satisfying

$$d_{L_p}(Q^0, Q) < \delta(\epsilon, n_0)$$

it is implied:

$$d_{L_p}(D^0(\hat{S}_n), D(\hat{S}_n)) < \epsilon ; \forall n \leq n_0$$

---

For data samples that are unlimited in number, we would like to investigate properties of the estimates that, in addition to satisfying the conditions in definition 4, they also guarantee fast convergence. Such an analysis will provide the designer with the additional valuable information of the sample sizes necessary to satisfy a given performance. The performance measure in

this case is, of course, the Lèvy deviation of the estimate whenever the data distributions move inside a certain sphere.

Before we proceed in the analysis, we need the assumption of a certain dependence structure of the data. As in the previous section, we will assume that the data vector  $X_n$  consists of  $n_1$   $k$ -dimensional independent vectors  $X_{kj}$ ;  $j = 1, \dots, n_1$ ;  $n = kn_1$ .

To avoid unnecessary generalities and to make the analysis more meaningful, we will also consider a particular, quite general form of estimates. Specifically, we will assume that to estimate the component  $s_i$  of the parameter vector  $S_p$  we apply a different, in general, continuous transformation on each of the first  $q$   $k$ -size data blocks and form a linear combination of these transformations. We repeat the same transformation to the next  $qk$  data block, etc. and we finally average out the resulting values. The assumption is, of course, that we always receive data in  $qk$  size blocks. To express the above description mathematically, we write:

$$\hat{s}_i(X_{nqk}) = \frac{1}{n} \sum_{j=1}^n \sum_{\ell=1}^q a_{\ell} \mu_{i\ell}(X_{k,j\ell}) \quad (27)$$

where  $X_{k,j\ell}$  the  $j\ell$ th  $k$ -size block of data from the vector  $X_{nqk}$ , and  $\mu_{i\ell}(\cdot)$  a continuous scalar function on the  $E^k$  Euclidean space for every  $1 \leq \ell \leq q$ . The continuity of  $\mu_{i\ell}(\cdot)$  guarantees satisfaction of lemma 4.

We consider estimators of the same general nature for all the  $S_p$  components, therefore we finally obtain a system of estimates as given by (27) for  $1 \leq i \leq p$ . We will observe at this point that since the  $k$ -dimensional data vectors  $X_{kj}$  have been assumed to be independent from each other, the functions  $\mu_{i\ell}(X_{k,j\ell})$  are independent random variables for different  $j\ell$ 's.

For convenience, we will pick here the distortion measure  $\rho$  that is included in the Lèvy distances of definition 4 to be given by:

$$\rho_0(X_n, Y_n) = \frac{1}{n} \sum_{i=1}^n |x_i - y_i|; X_n = \{x_i; i = 1, \dots, n\};$$

$$Y_n = \{y_i; i = 1, \dots, n\} \quad (28)$$

Some different distortion (penalty) measures  $\rho(\cdot, \cdot)$  lead to analysis similar to the one that will be presented in this section, therefore the choice in (28) is not truly restrictive.

Starting on the analysis of the estimators in (27), let us suppose that for some  $\epsilon > 0$ , some  $n$ , and some  $Q^0, Q^1$  distributions we have:

$$d_{L_{\rho_0}}(D^0(\hat{S}_{nqk}), D^1(\hat{S}_{nqk})) < \epsilon \quad (29)$$

According to Strassen ([2], Th. 11) the condition (29) is true if and only if there is some  $2p$ -dimensional distribution  $D^{01}(\cdot, \cdot)$  with  $D^0(\cdot), D^1(\cdot)$  marginals such that

$$D^{01}\left\{\frac{1}{p} \sum_{i=1}^p \frac{1}{n} \left| \sum_{j=1}^n \sum_{l=1}^q a_l [\mu_{il}(X_{k,jl}) - \mu_{il}(Y_{k,jl})] \right| \geq \epsilon\right\} \leq \epsilon \quad (30)$$

In (30), the estimate form in (27) is considered and  $X_{k,jl}, Y_{k,jl}$  are distributed as in  $Q^0, Q^1$  respectively.

The joint distribution  $D^{01}$  can be translated to a distribution  $Q^{01}$  with  $Q^0, Q^1$  marginals instead through a specific estimate choice. Since the  $Q^0, Q^1$  distributions are representing data of  $k$ -size independent blocks,  $Q^{01}$  and  $D^{01}$  will be such that they maintain this independence. In other words,  $D^{01}$  in (30) should be such that the differences  $[\mu_{il}(X_{k,jl}) - \mu_{il}(Y_{k,jl})]$  are independent from each other for different  $jl$  values.

Observing expression (30), we see that due to the truth of the inequality:

$$\sum_{i=1}^p D^{01}\left\{\frac{1}{n} \sum_{j=1}^n \sum_{l=1}^q a_l [\mu_{il}(X_{k,jl}) - \mu_{il}(Y_{k,jl})] \geq \epsilon\right\} \geq$$

$$\geq D^{ol}\left\{\frac{1}{n} \left| \sum_{j=1}^n \sum_{l=1}^q a_l [\mu_{il}(X_{k,jl}) - \mu_{il}(Y_{k,jl})] \right| \geq \epsilon \right\}$$

$$\geq D^{ol}\left\{\frac{1}{p} \sum_{i=1}^p \frac{1}{n} \left| \sum_{j=1}^n \sum_{l=1}^q a_l [\mu_{il}(X_{k,jl}) - \mu_{il}(Y_{k,jl})] \right| \geq \epsilon \right\}$$

to have (29) satisfied, it is sufficient that:

$$D^{ol}\left\{\frac{1}{n} \left| \sum_j \sum_{l=1}^q a_l [\mu_{il}(X_{k,jl}) - \mu_{il}(Y_{k,jl})] \right| \geq \epsilon \right\} \leq \frac{\epsilon}{p} ; \forall i: 1 \leq i \leq p \quad (31)$$

So, if for some  $D^{ol}$  choice with  $Q^0, Q$  marginals, (31) is satisfied, so is (29). For additional simplification we will assume that we are working on distribution spaces that generate stationary data. Then,

$$Q^0(X_{k,jl} = Z_k) = Q^0(X_{k,mr} = Z_k) ; \forall jl, mr$$

$$Q^1(Y_{k,jl} = Z_k) = Q^1(Y_{k,mr} = Z_k) ; \forall jl, mr$$

In this case the distances  $d_L(Q^0, Q), d_L(D^0(\hat{S}_n), D(\hat{S}_n))$  are generated by the  $k$ -dimensional distributions  $Q_k^0, Q_k^{1p}$ . Also, we can then define:

$$m_{il}^0 = E_{Q^0}\{\mu_{il}(X_{k,jl})\}$$

$$m_{il}^1 = E_{Q^1}\{\mu_{il}(Y_{k,jl})\} \quad (32)$$

where  $m_{il}^0, m_{il}^1$  are independent of  $jl$ . From (32) and (31) we obtain that if we want (29) satisfied, it is sufficient to require the following condition:

$$D^{ol}\left\{\left| \frac{1}{n} \sum_{j=1}^n \sum_{l=1}^q a_l [\mu_{il}(X_{k,jl}) - m_{il}^0] - [\mu_{il}(Y_{k,jl}) - m_{il}^1] \right| + \right.$$

$$\left. + \sum_{l=1}^q a_l [m_{il}^0 - m_{il}^1] \right| \geq \epsilon \right\} \leq \frac{\epsilon}{p} ; \forall i: 1 \leq i \leq p \quad (33)$$

where

$$E_{D^{ol}}\{\mu_{il}(X_{k,jl}) - \mu_{il}(Y_{k,jl})\} = m_{il}^0 - m_{il}^1 \quad (34)$$

Directly from (33) we can express the following stronger condition, to guarantee

satisfaction of (29):

$$D_k^{ol} \left\{ \left| \frac{1}{n} \sum_{j=1}^n \sum_{\ell=1}^q a_{\ell} [\mu_{i\ell}(X_{k,j\ell}) - m_{i\ell}^0] - [\mu_{i\ell}(Y_{k,j\ell}) - m_{i\ell}^1] \right| \right. \\ \left. \geq \epsilon - \sum_{\ell=1}^q |a_{\ell}| |m_{i\ell}^0 - m_{i\ell}^1| \right\} \leq \frac{\epsilon}{p}; \quad \forall i: 1 \leq i \leq p \quad (35)$$

At this point we find necessary to summarize our analysis up to this point in the following Corollary:

#### Corollary 1

If, given  $\epsilon > 0$ , there is some  $n_a$  and some  $\delta(\epsilon, n_a) > 0$  such that:  
 $\forall n \geq n_a$  and every  $Q_k^1$  satisfying:

$$d_{L_{\rho_0}}(Q_k^0, Q_k^1) < \delta(\epsilon, n_a)$$

there is some  $D_k^{ol}$  with marginals  $Q_k^0, Q_k^1$  implying:

$$D_k^{ol} \left\{ \left| \frac{1}{n} \sum_{j=1}^n \sum_{\ell=1}^q a_{\ell} [\mu_{i\ell}(X_{k,j\ell}) - m_{i\ell}^0] - [\mu_{i\ell}(Y_{k,j\ell}) - m_{i\ell}^1] \right| \right. \\ \left. \geq \epsilon - \sum_{\ell=1}^q |a_{\ell}| |m_{i\ell}^0 - m_{i\ell}^1| \right\} \leq \frac{\epsilon}{p}; \quad \forall i: 1 \leq i \leq p \quad (36)$$

where  $\mu_{i\ell}(\cdot)$  continuous on  $E^k$  for  $1 \leq \ell \leq q; 1 \leq i \leq p$ , then the estimate described by (27) is Lèvy-Lèvy weakly robust according to definition 4.

---

Now that we have summarized the observations up to now, we will continue with analysis of the conditions of corollary 1.

Observing condition (36) we realize that if we want it satisfied for given  $\epsilon > 0$ , the sum

$$\sum_{\ell=1}^q |a_{\ell}| |m_{i\ell}^0 - m_{i\ell}^1| \quad (37)$$



must be smaller than  $\epsilon$ . In other words, for given  $\epsilon > 0$ ,  $n_a$  we would like to find some  $\delta(\epsilon, n_a) > 0$  that first guarantees the satisfaction of the inequality

$$\sum_{\ell=1}^q |a_{\ell}| |m_{i\ell}^0 - m_{i\ell}^1| < \epsilon_1 ; \text{ where } \epsilon_1 > 0 \text{ and such that } \epsilon_1 < \epsilon \quad (38)$$

for every  $Q_k^1: d_{L,p}(Q_k^0, Q_k^1) < \delta(\epsilon, n_a)$  and then it secures the existence of some  $D_k^{01}$  with  $Q_k^0, Q_k^1$  marginals that satisfies (36).

Let us pick  $\epsilon_1 = \epsilon/2$ . Then, if  $\mu_{i\ell}(X_k)$  is bounded absolutely by some constant  $B$  for every  $1 \leq i \leq p, 1 \leq \ell \leq q$  and  $X_k$ , we know from lemma 3 and its proof in appendix A that given

$$\frac{\epsilon}{2 \sum_{\ell=1}^q |a_{\ell}|}$$

there is some

$$\delta\left(\frac{\epsilon}{2 \sum_{\ell=1}^q |a_{\ell}|}, k, B\right) > 0$$

such that

$$\sum_{\ell=1}^q |a_{\ell}| |m_{i\ell}^0 - m_{i\ell}^1| < \frac{\epsilon}{2}$$

After this last observation, we can go one step further and express the following corollary that is a simplified form of corollary 1:

### Corollary 2

Let the set of  $\mu_{i\ell}(X_k)$  estimates in (27) be continuous on  $E^k$  and absolutely bounded by some  $B > 0$  for every  $1 \leq \ell \leq q, 1 \leq i \leq p$ .

Let

$$\delta\left(\frac{\epsilon}{\frac{q}{2 \sum_{\ell=1}^q |a_{\ell}|}}, k, B\right) > 0$$

be a constant such that: if  $Q_k^1$  is satisfying

$$d_{L, \rho_0}(Q_k^0, Q_k^1) < \delta\left(\frac{\epsilon}{\frac{q}{2 \sum_{\ell=1}^q |a_{\ell}|}}, k, B\right)$$

then

$$\sum_{\ell=1}^q |a_{\ell}| |m_{i\ell}^0 - m_{i\ell}^1| < \frac{\epsilon}{2}$$

for some given  $\epsilon > 0$ . In this case, if for the same above given  $\epsilon > 0$ , and for some natural number  $n_a$ , there is some  $\delta(\epsilon, n_a) > 0$ :

$$\delta(\epsilon, n_a) \leq \delta\left(\frac{\epsilon}{\frac{q}{2 \sum_{\ell=1}^q |a_{\ell}|}}, k, B\right)$$

that for every  $n \geq n_a$  and  $Q_k^1$  satisfying

$$d_{L, \rho_0}(Q_k^0, Q_k^1) < \delta(\epsilon, n_a)$$

there is some  $D_k^{01}$  with  $Q_k^0, Q_k^1$  marginals that satisfies the condition:

$$D_k^{01}\left\{\left|\frac{1}{n} \sum_{j=1}^n \sum_{\ell=1}^q a_{\ell} [\mu_{i\ell}(X_{k,j\ell}) - m_{i\ell}^0] - [\mu_{i\ell}(Y_{k,j\ell}) - m_{i\ell}^1]\right| \geq \frac{\epsilon}{2}\right\} \leq \frac{\epsilon}{p}$$

$$\forall 1 \leq i \leq p \quad (39)$$

Then the estimate in (27) is Lévy-Lévy weakly robust by definition 4.

---

The boundness condition on the estimates that appeared in section 4 is included again in corollary 2. As mentioned before, this condition is a realistic property that protects the estimate values from wandering in the

space of unacceptable values.

We will now concentrate our attention on expression (39), hoping to obtain a comparatively small lower bound on the sample size  $n$  that satisfies the bound  $\epsilon/p$ . That, of course, we expect for particular  $\mu_{il}(\cdot)$  choices.

We will point out again that the brackets

$$[[\mu_{il}(X_{k,jl}) - m_{il}^0] - [\mu_{il}(Y_{k,jl}) - m_{il}^1]] \quad (40)$$

in (39) are zero mean independent variables for every different  $jl$  value.

A theorem expressed by Revész ([5], pg. 57) will be extremely useful here. We state the theorem below.

#### Theorem 1

Let  $x_1, x_2, \dots, x_n$  be independent, zero mean, not necessarily identically distributed variables. Then, the probability

$$P_n(\eta) = P_x \left\{ \left| \frac{x_1 + x_2 + \dots + x_n}{n} \right| \geq \eta \right\}$$

converges to zero exponentially for any  $\eta > 0$ , i.e., there is some  $C > 0$  and some  $0 < \nu < 1$  such that:

$$P_n(\eta) \leq C\nu^n$$

if and only if: For all  $\eta > 0$  there exists a constant  $\epsilon_\eta > 0$  and some  $t_\eta > 0$  such that:

$$\prod_{k=1}^n E\{e^{tx_k}\} \leq \epsilon_\eta e^{|t|\eta n} \quad \text{whenever} \quad |t| \leq t_\eta$$

Also, the probability  $P_n(\eta)$  cannot converge faster than exponentially and the constants  $C$  and  $\nu$  that express the bound on the probability  $P_n(\eta)$  are chosen as follows:  $C = \epsilon_\eta$ ;  $\nu = e^{-(\delta-\eta) \cdot t_\eta}$ ; where  $\delta$  some value in the interval  $(0, \eta)$ .

We will apply the above theorem on condition (39) to design robust estimates that, in addition to satisfying the continuity and boundness properties expressed in corollary 2, they also have the strong characteristic of the fastest possible convergence to their asymptotic value, which through the absolute boundness is guaranteed to be stable inside stationary  $Q_k^0$ -Lévy-contaminated distribution families, whenever the contamination is small enough.

Directly from theorem 1, from the fact that the expressions in (40), are independent and zero mean for different  $j\ell$  values and from the observation that the sums

$$\sum_{\ell=1}^q a_{\ell} [\mu_{i\ell}(X_{k,j\ell}) - m_{i\ell}^0] - [\mu_{i\ell}(Y_{k,j\ell}) - m_{i\ell}^1] \quad (41)$$

are identically distributed for every  $j$ , we obtain that the left part of inequality (39) converges exponentially if and only if:

For all  $\epsilon > 0$  there exists some constant  $\epsilon_{\epsilon/2} > 0$  and some  $t_{\epsilon/2} > 0$  such that

$$\left[ \prod_{1 \leq \ell \leq q} E_{D_k^{ol}} \left\{ e^{t a_{\ell} [\mu_{i\ell}(X_k) - m_{i\ell}^0] - [\mu_{i\ell}(Y_k) - m_{i\ell}^1]} \right\} \right]^n \leq \epsilon_{\epsilon/2} e^{|t| \frac{\epsilon}{2} n} \quad (42)$$

$$\forall |t| \leq t_{\epsilon/2}$$

Due to theorem 1, the larger  $t_{\epsilon/2}$  we can find the faster the convergence of the estimators to their asymptotic value. Also, we must emphasize here that we are seeking a  $t_{\epsilon/2}$  that is common for all  $Q_k^1$  that are members of the data distribution contaminated family.

As an additional observation on (42), we see that since its left part is equal to one for  $t = 0$ ,  $\epsilon_{\epsilon/2}$  cannot be smaller than one. Seeking the smallest possible  $\epsilon_{\epsilon/2}$  we may as well pick  $\epsilon_{\epsilon/2} = 1$ .

In this case, condition (42) can take the form:

$$\sum_{\ell=1}^q \ln E_{D_k^{ol}} \{ e^{t a_{\ell} [\mu_{i\ell}(X_k) - m_{i\ell}^0] - [\mu_{i\ell}(Y_k) - m_{i\ell}^1]} \} \leq |t| \frac{\epsilon}{2} \quad \forall |t| \leq t_{\epsilon/2} \quad (43)$$

where  $t_{\epsilon/2}$  some positive value and  $D_k^{ol}$  some  $2k$ -dimensional distribution with  $Q_k^0, Q_k^1$  marginals. Each of the logarithmic expressions in (43) is a convex  $U$  function of  $t$  with minimum at  $t = 0$  and minimum value equal to zero. This is true due to the fact that each of these logarithmic functions has positive second derivative for every  $t$  and first derivative at  $t = 0$  that is equal to zero. The above observations are true for every  $Q_k^0, Q_k^1, D_k^{ol}$  choice. Also, the sum of the logarithmic functions in (43) is also convex  $U$  with minimum that is equal to zero and happens at  $t = 0$ .

Due to the above observations, there is always some  $t_{\epsilon/2}$  (for given  $\epsilon > 0$ ) that satisfies (43). The analysis should be concentrated now to designing the constants  $a_{\ell}$  and the functions  $\mu_{i\ell}(\cdot)$  in a way that will make the common for all distributions in the contaminated family  $t_{\epsilon/2}$  as large as possible.

Before we express some thoughts on that we will state the following theorem:

### Theorem 2

Let the set of estimates in (27) be continuous  $E^k$  and absolutely bounded by some  $B > 0$  for every  $1 \leq \ell \leq q; 1 \leq i \leq p$ . Then, if for some  $\delta > 0$  and such that it is nonlarger than the

$$\delta \left( \frac{\epsilon}{2 \sum_{\ell=1}^q |a_{\ell}|}, k, B \right)$$

included in corollary 2, a common  $t_{\epsilon/2} > 0$  can be found for all the members of the  $Q_k^0$ -contaminated family characterized by:

$$D_{L\rho_0}(Q_k^0, Q_k^1) \leq \delta \quad (44)$$

where  $t_{\epsilon/2}$  is such that it satisfies the expression

$$\sum_{l=1}^q \ln E_{Q_k^0} \{ e^{t_{\epsilon/2} [\mu_{il}(X_k) - m_{il}^0]} \} E_{Q_k^1} \{ e^{-t_{\epsilon/2} [\mu_{il}(X_k) - m_{il}^1]} \} \leq |t| \frac{\epsilon}{2} \quad (45)$$

$$\forall |t| \leq t_{\epsilon/2}; \forall Q_k^1 \text{ satisfying (44)}$$

the estimate in (27) is weakly robust and it converges to its asymptotic value with rate expressed by the bound:

$$e^{n(\zeta - \frac{\epsilon}{2})t_{\epsilon/2}}; \zeta \in (0, \frac{\epsilon}{2}) \quad (46)$$

The absolute boundness of  $\mu_{il}(\cdot)$  guarantees stability of the  $\mu_{il}(\cdot)$  asymptotic value inside the contaminated family expressed by (44).

Also, for performance equal to  $\epsilon$  (where  $\epsilon > 0$ ) it is sufficient that we choose then number of samples equal to:

$$n_a = \frac{\ln \frac{\epsilon}{p}}{(\zeta - \frac{\epsilon}{2})t_{\epsilon/2}} \quad (47)$$

$\zeta$  is some value in  $(0, \epsilon/2)$  and it is the same in both expressions (46) and (47).

---

In theorem 2 we picked  $D_k^{01}(X, Y) = Q_k^0(X) \cdot Q_k^1(Y)$ . We will now concentrate on expression (45) to draw some useful conclusions. Indeed we can separate the expression (45) into two parts: One including the behavior of the estimate at the central distribution  $Q_k^0$  and one describing the same behavior at the distributions included in the contaminated family in (44).

Through the separation we just mentioned we can write (45) in the following way:

$$\sum_{l=1}^q \ln E_{Q_k^0} \{ e^{t a_l [\mu_{il}(X_k) - m_{il}^0]} \} + \sum_{l=1}^q \ln E_{Q_k^1} \{ e^{-t a_l [\mu_{il}(X_k) - m_{il}^1]} \} \leq |t| \frac{\epsilon}{2}$$

$$\forall |t| \leq t_{\epsilon/2} \quad (48)$$

Both sums on the left part of (48) are convex U with minimum equal to zero and assumed at  $t = 0$ . The less sharp both of these sums are as functions of  $t$ , the larger  $t_{\epsilon/2}$  will be for given  $\epsilon > 0$ . Since  $Q_k^0$  is a well-defined distribution, the set  $\{a_l, \mu_{il}(\cdot); 1 \leq l \leq q, 1 \leq i \leq p\}$  can be designed to make the sum

$$\sum_{l=1}^q \ln E_{Q_k^0} \{ e^{t a_l [\mu_{il}(X_k) - m_{il}^0]} \} \quad (49)$$

as less sharp at  $t = 0$  as possible.

We will call an estimate  $\{a_l, \mu_{il}(X_k); 1 \leq l \leq q; 1 \leq i \leq p\}$  in (27) a fast estimate at  $Q_k^0$  if it makes the expression in (49) a slowly increasing function of  $t$  around  $t = 0$ . The continuity of the functions  $\mu_{il}(\cdot)$  on  $E^k$  and their boundness, guarantees closeness of the moment generating functions

$$E_{Q_k^1} \{ e^{-t a_l [\mu_{il}(X_k) - m_{il}^0]} \}$$

to the central moment generating function  $E_{Q_k^0} \{ e^{-t a_l [\mu_{il}(X_k) - m_{il}^0]} \}$

for  $Q_k^1$  belonging to the contaminated family (44). Therefore, for fast ex-

ponential convergence, then, it is sufficient to design the estimator in (27)

to make the logarithmic expression of the moment generating functions at

$Q_k^0$  (expression (49)) as slowly increasing with  $t$  around  $t = 0$  as possible.

As conclusion, we finally express the following theorem:

### Theorem 3

If the set of estimates in (27) is continuous as a function on  $E^k$ , absolutely bounded for every  $1 \leq l \leq q, 1 \leq i \leq p$  and it is a fast estimate

at  $Q_k^0$ , then it is also weakly robust at  $Q_k^0$ , and it converges fast and exponentially to its stable mean.

The advantage of the method presented in this section, in comparison to Hampel's method is that the design of the weakly robust estimates reduces to real function properties and to moment generating functions at the central distribution  $Q_k^0$ . Furthermore, this method has the tremendous advantage of treating finite data samples (through the smallest  $n_a$  that achieves the required performance) and not just asymptotic situations as in [1]. Finally, the Lèvy-performance criterion is stronger than the small-variance criterion treated by Huber [1].

Finding the smallest  $t_{e/2}$  in theorem 2 that will satisfy all members of the  $Q_k^0$ -contaminated family is a task, samples of which will be shown in the following section.

Concluding this section, we will mention that theorem 1 has also been applied in [19] to find confidence intervals for Bayes error estimation which were subsequently used to determine the optimal degree of quantization.

## 6. FURTHER DISCUSSION ON SECTIONS 3, 4 AND 5

In this section we will present some discussion on the application of the theory that appears in the previous sections. Our discussion will be mostly oriented toward expressing methods and suggesting ideas for the actual design of the robust estimators under consideration. The possibilities for such designs are many and lengthy and they will appear in detail in future work (under preparation).

1. We will first start with the Vasershtein robust estimators that are analyzed in section 3. According to lemma 2, such an estimator (where the underlying distortion measure is the square error) must be linear. That is:



$$\hat{S}_n(X_n) = \{\hat{s}_{in}(X_n); i = 1, \dots, p\}; \hat{s}_{in}(X_n) = \sum_{k=1}^n a_{ni}(k)x_k \quad (50)$$

Also, a sufficient condition for robustness in this case is that each of the sequences

$$\left\{ \sum_{k=1}^n |a_{ni}(k)| \right\}; i = 1, \dots, p$$

converges to some finite value. Any

$$\left\{ \sum_{k=1}^n |a_{ni}(k)| \right\}$$

that is geometric with multiplying factor  $\omega$  smaller than one belongs to the above group. Also, the experimental mean is obviously Vasershtein robust.

In general, one will look for linear estimates that converge at the central well-known distribution to the desired value and whose coefficients satisfy the absolute convergence property mentioned above. For example, let the central distribution  $Q^0$  be  $m$ -dependent. Then accumulate the data in  $m$ -groups and form the estimate

$$\hat{s}_{i, nm}(X_{nm}) = \frac{1}{n} \sum_{k=1}^{nm} a_{nm,i}^1(k)x_k \quad (51)$$

where the previous coefficients  $a_{ni}(k)$  in (50) are related to the  $a_{nm,i}^1(k)$ 's in (51) by

$$a_{nm,i}(k) = \frac{1}{n} a_{nm,i}^1(k) \quad (52)$$

Due to the  $m$ -dependence of  $Q^0$ , the sums

$$\sum_{k=1}^m a_{nm,i}^1(k)x_k, \dots, \sum_{k=jm}^{(j+1)m} a_{nm,i}^1(k)x_k, \dots$$

are independent from each other and according to theorem 1, the estimate in (51) converges exponentially and in probability to

$$m_0 \cdot \lim_{n \rightarrow \infty} n^{-1} \sum_{k=1}^n a_{nm,i}^{(k)} = m_0 \cdot \lim_{n \rightarrow \infty} \sum_{k=1}^n a_{ni}^{(k)} \quad (53)$$

where  $m_0$  the one-dimensional mean determined by  $Q^0$ .

If we pick all the coefficients positive and design the sum  $\sum_{k=1}^n a_{ni}^{(k)}$  so that it converges to the desired value at  $Q^0$ , then at the same time the robustness sufficient condition is also satisfied and the estimate in (51) is stable inside the Vasershtein contaminated family.

2. For the Lèvy-Vasershtein robustness (as found in section 4), one should seek for continuous as real functions and continuous at the central distribution estimates that are also absolutely bounded. This last property is the only additional to the ones required by the robust in the Hampel sense estimates.

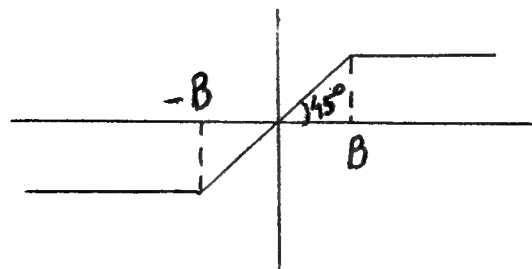
Of course, all estimates that are not Hampel robust (such as the experimental mean), are not robust in the present sense either.

Also, the L-estimators

$$\hat{s}_{in}(X_n) = \sum_{j=1}^n b_j x_{(j)} \quad (54)$$

where  $X_n = \{x_i\}$ ,  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ , the ordered components of  $X_n$ , are not in general absolutely bounded, therefore not Lèvy-Vasershtein robust in general.

On the contrary, a truncated version of the estimates in (54) can be Lèvy-Vasershtein robust with proper choice of the coefficients  $b_j$ . By truncated version of (54) we mean that we substitute each  $x_{(j)}$  in it by its output when passed through the following nonlinearity:



Finally, using the approach of the M-estimators, we can find the estimate  $\hat{s}_{in}(X_n)$  as the truncated zeros of a sum [18]:

$$\sum_{j=1}^n \psi_1(x_j - \hat{s}_{in}(X_n)) \quad (55)$$

where  $\psi_1(\cdot)$  a smooth enough function. The zeros of the expression in (55) will be passed through a nonlinearity as the one mentioned above, or any other that cuts off all estimate values that are absolutely higher than B.

3. The fast robustness introduced in section 5 is important and deserves special attention. Indeed, the method introduced there can be considered small sample since the effort is toward designing estimates that, besides being robust, converge to their stable (for small deviations of the data distribution from the central one) mean fast.

The measure of performance is Lèvy distance. In other words, if a specific contaminated family of data distributions is given, the stability of the estimate in it is measured through the longest Lèvy distance of the estimate distribution at some  $Q$  from the estimate distribution at  $Q^0$ , where  $Q^0$  is the central well-known data distribution and  $Q$  moves inside the contaminated family.

As explained in section 5, if  $p$  scalar estimates are calculated from the same data and the performance required is a fixed  $\epsilon > 0$ , then the minimum number of data  $n_a$  that will satisfy this performance is given by

$$n_a = \frac{\ln \frac{\epsilon}{p}}{(\zeta - \frac{\epsilon}{2}) t_{\epsilon/4}} \quad (56)$$

where  $\zeta$  is some arbitrary constant in  $(0, \epsilon/2)$  and  $t_{\epsilon/4}$  is for given estimates (of (27) form and  $k$ -dependent group data) the smallest  $t_{Q_k^1, \epsilon/4}$  such that

$$\sum_{\ell=1}^q \ln E_{Q_k^1} \{ e^{-t a_{\ell} [\mu_{i\ell}(X_k) - m_{i\ell}^1]} \} \leq |t| \frac{\epsilon}{4} ; \forall |t| \leq t_{Q_k^1, \epsilon/4} \quad (57)$$

where  $Q_k^1$  moves in the considered data contaminated family and the smallest  $t_{Q_k^1, \epsilon/4}$  is taken among these  $Q_k^1$  members.

To achieve this minimum  $t_{\epsilon/4}$  as large as possible, we go backward. Specifically, we are looking for this estimate set  $\{a_{\ell}; \mu_{i\ell}(X_k); 1 \leq \ell \leq q; 1 \leq i \leq p\}$  that achieves the largest  $t_{Q_k^1, \epsilon/4}$  for the worst  $Q_k^1$  choice inside the specific given contaminated family. This worst case corresponds to

$$\sum_{\ell=1}^q \ln E_{a_k} \{ e^{-t a_{\ell} [\mu_{i\ell}(X_k) - m_{i\ell}^1]} \}$$

function with the fastest increase for  $|t|$  increasing.

Specific continuous and absolutely bounded functions  $\mu_{i\ell}(X_k)$  can be chosen with some of their characteristics left flexible. The design of these characteristics becomes a maximum problem (find largest  $t_{Q_k^1, \epsilon/4}$  for the

logarithmic moment generating expression the sharpest in the family) and the methods for solving it are similar to the ones applied in [1] and [19].

# REFERENCES

1. P.J. Huber (1964) "Robust estimation of a location parameter," Ann. Math. Statist. 36: 1753-1758.
2. V. Strassen (1965) "The existence of probability measures with given marginals," Ann. Math. Statist. 36: 423-439.
3. P.J. Huber (1967) "The behavior of maximum likelihood estimates under nonstandard conditions," Proc. Fifth Berkeley Symp. Math. Statist. Prob. 1: 221-233.
4. R.M. Dudley (1968) "Distances of probability measures and random variables," Ann. Math. Statist. 39: 1563-1572.
5. Fál Revész (1968) The Laws of Large Numbers in Probability and Mathematical Statistics. A series of monographs and textbooks edited by Z.W. Birnbaum and E. Lukacs, Academic Press, New York.
6. L.N. Vasershtein (1969) "Markov processes on a countable product space, describing large systems of automata," Problemy Peredachi Informatsii, 5(3): 64-73 (in Russian).
7. P.J. Huber (1970) "Studentizing Robust Estimates," in Nonparametric Techniques in Statistical Inference, M.L. Puri, ed., Cambridge Univ. Press: 453-463.
8. R.L. Dobrushin (1970) "Prescribing a system of random variables by conditional distributions," Theory of Prob. and Appl. XV(3): 458-486.
9. L.N. Vasershtein and A.M. Leonfovich (1970) "On invariant measures of certain Markov operators, describing a random medium," Problemy Peredachi Informatsii 6(1): 71-80 (in Russian).
10. F.R. Hampel (1971) "A general qualitative definition of robustness," Ann. Math. Statist. 42: 1887-1896.
11. R.D. Martin and S.C. Schwartz (1971) "Robust detection of a known signal in nearly Gaussian noise," IEEE Trans. Inf. Th. IT-17: 50-56.

12. P.J. Huber (1972) "Robust statistics. A review," Ann. Math. Statist. 43: 1041-1067.
13. R.D. Martin and C.J. Masreliez (1975) "Robust estimation via stochastic approximation," IEEE Trans. Inf. Th. IT-21: 257-262.
14. R.M. Gray, D. Neuhoff and P.C. Shields (1975) "A generalization of Orstein's distance with applications to information theory," Ann. Prob. 3(2): 315-324.
15. D.L. Neuhoff, R.M. Gray and L.D. Davisson (1975) "Fixed rate universal block source coding with a fidelity criterion," IEEE Trans. Inf. Th. IT-21(5): 511-523.
16. P. Papantoni-Kazakos (1975) "Small sample efficiencies of rank tests," IEEE Trans. Inf. Th. IT-21: 150-156.
17. P. Papantoni-Kazakos (1976) "Some distance measures and their use in feature selection," Rice University Technical Report #7611, Electrical Engineering Department, Rice University, Houston, Texas.
18. P. Papantoni-Kazakos (1977) "Robustness in parameter estimation," IEEE Trans. Inf. Th., accepted for publication, March 1977.
19. D. Kazakos (1977) "Quantization complexity and training data sample size in detection problems," submitted for publication.
20. S.A. Kassam and J.B. Thomas (1975) "A class of nonparametric detectors of dependent input data," IEEE Trans. Inf. Th. IT-21: 431-437.
21. S.A. Kassam and J.B. Thomas (1976) "Asymptotically robust detection of a known signal in contaminated non-Gaussian noise," IEEE Trans. Inf. Th. IT-22: 22-26.

APPENDIXProof of Lemma 3

In the distance  $d_{Vt}(D^0(\hat{S}_n), D(\hat{S}))$  in (25), the terms

$$(\sqrt{\sigma_{in}^0(0)} - \sqrt{\sigma_{in}(0)})^2$$

and  $(m_{in}^0 - m_{in})^2$  appear. We develop upper bounds on each of them. Indeed, we obtain:

$$(\sqrt{\sigma_{in}^0(0)} - \sqrt{\sigma_{in}(0)})^2 < |\sqrt{\sigma_{in}^0(0)} - \sqrt{\sigma_{in}(0)}| |\sqrt{\sigma_{in}^0(0)} + \sqrt{\sigma_{in}(0)}|$$

$$\begin{aligned} &= |\sigma_{in}^0(0) - \sigma_{in}(0)| = \left| \int \hat{s}_{in}^2(X_n) Q_0(dX_n) - \int \hat{s}_{in}^2(X_n) Q(dX_n) \right. \\ &\quad \left. - [(m_{in}^0)^2 - (m_{in})^2] \right| \\ &\leq \left| \int \hat{s}_{in}^2(X_n) Q_0(dX_n) - \int \hat{s}_{in}^2(X_n) Q(dX_n) \right| + |(m_{in}^0)^2 - (m_{in})^2| \quad (A.1) \end{aligned}$$

If the estimate  $\hat{s}_{in}(X_n)$  is absolutely bounded from above by some  $A_1 > 0$ , then we obtain from (A.1):

$$\begin{aligned} (\sqrt{\sigma_{in}^0(0)} - \sqrt{\sigma_{in}(0)})^2 &< \int |\hat{s}_{in}(X_n) + \hat{s}_{in}(Y_n)| |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| \cdot \\ &\cdot D(dX_n, dY_n) + 2A_1 \int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \\ &\leq 4A_1 \int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \quad (A.2) \end{aligned}$$

where  $D(X_n, Y_n)$  some distribution with  $Q^0(X_n)$ ,  $Q(Y_n)$  marginals. From (A.2)

we also obtain, for the same  $D(X_n, Y_n)$ :

$$\begin{aligned} &(\sqrt{\sigma_{in}^0(0)} - \sqrt{\sigma_{in}(0)})^2 + (m_{in}^0 - m_{in})^2 \\ &\leq 6A_1 \int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \end{aligned}$$

For robustness from  $d_L(Q_0, Q)$  to  $d_V(D^0(\hat{S}_n), D(\hat{S}_n))$  we require:

That given  $\epsilon > 0$ , there is some  $\delta(\epsilon) > 0$  such that

If  $d_L(Q_0, Q) < \delta(\epsilon) \Rightarrow d_V(D^0(\hat{S}_n), D(\hat{S}_n)) < \epsilon ; \forall n$ .

Or:

Given  $\epsilon > 0$ , there is some  $\delta(\epsilon) > 0$ :

For  $d_L(Q_0, Q) < \delta(\epsilon)$

$$\Rightarrow \max_{1 \leq i \leq p} (6A_i \int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n)) < \epsilon; \forall n$$

where  $D(X_n, Y_n)$  has respective marginals  $Q_n^0(X_n)$  and  $Q_n(Y_n)$ . If given

$\epsilon > 0$ , there is for every  $1 \leq i \leq p$  some  $\delta_i(\epsilon) > 0$ :

$$\int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) < \frac{\epsilon}{6A_i}; \forall n$$

Then

$$\delta(\epsilon) = \min_{1 \leq i \leq p} \delta_i(\epsilon)$$

If  $d_L(Q_n^0, Q_n) < \delta_{in}(\epsilon)$ , then there is some  $\forall n$   $D(X_n, Y_n)$  with  $Q_n^0(X_n), Q_n(Y_n)$  marginals (Strassen [2], Th. 11) such that

$$D(X_n, Y_n; \rho(X_n, Y_n)) \geq \delta_{in}(\epsilon)$$

Pick this  $D(X_n, Y_n)$  and write:

$$\begin{aligned} & \int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \\ &= \int_{X_n, Y_n: \rho(X_n, Y_n) < \delta_{in}(\epsilon)} |\hat{s}_{in}(Y_n) - \hat{s}_{in}(X_n)| D(dX_n, dY_n) + \int_{X_n, Y_n: \rho(X_n, Y_n) \geq \delta_{in}(\epsilon)} |\hat{s}_{in}(Y_n) - \hat{s}_{in}(X_n)| D(dX_n, dY_n) \\ &\leq 2A_i \delta_{in}(\epsilon) + \int_{X_n, Y_n: \rho(X_n, Y_n) < \delta_{in}(\epsilon)} |\hat{s}_{in}(Y_n) - \hat{s}_{in}(X_n)| D(dX_n, dY_n) \end{aligned}$$

since  $|\hat{s}_{in}(X_n)| \leq A_i; \forall n$

Let  $\hat{s}_{in}(X_n)$  be continuous function everywhere on  $E^n$ . Then, given

$\epsilon_i > 0$  and  $X_n$ , there is some  $\delta(\epsilon_i, X_n) > 0$  such that

$$\forall Y_n: \rho(X_n, Y_n) < \delta(\epsilon_i, X_n) \Rightarrow |\hat{s}_{in}(Y_n) - \hat{s}_{in}(X_n)| < \epsilon_i$$



Choose a sequence  $\{c_j\}$  of positive, monotonically decreasing toward zero components and define

$$A_j = \{X_n : \delta(\epsilon_j, X_n) > c_j\}$$

Then

$$A_j \subseteq A_{j+1} \rightarrow E^n ; UA_j = E^n$$

Define

$$U(X_n) = \{Y_n : \rho(X_n, Y_n) < \delta(\epsilon_1, X_n)\}$$

$$B_j = U \left\{ U(X_n) : X_n \in A_j \right\}$$

Then  $B_j \subseteq B_{j+1}$  ,  $UB_j = E^n$

Now, given  $\eta_1 > 0$  , there is some  $k_1(n, \eta_1)$  such that

$$Q_0 \left( \begin{matrix} UB_j \\ 1 \leq j \leq k_1(n, \eta_1) \end{matrix} \right) > 1 - \eta_1$$

Denote:  $\mathcal{E}_{in} = \begin{matrix} UB_j \\ 1 \leq j \leq k_1(n, \eta_1) \end{matrix}$

Pick

$$\epsilon_1 = \frac{\epsilon}{18A_1} , \quad \eta_1 = \frac{\epsilon}{36A_1^2}$$

$$\delta_{in}(\epsilon) = \min \left\{ c_{k_1} \left( n, \frac{\epsilon}{36A_1^2} \right), \frac{\epsilon}{36A_1^2} \right\}$$

Then

$$\begin{aligned}
& \int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \\
& \leq 2A_i \delta_{in}(\epsilon) + \int_{\substack{X_n \in \mathcal{E}_{in} \\ Y_n: \rho(X_n, Y_n) < \delta_{in}(\epsilon)}} |\hat{s}_{in}(Y_n) - \hat{s}_{in}(X_n)| D(dX_n, dY_n) \\
& + \int_{\substack{X_n \in \mathcal{E}_{in}^c \\ Y_n: \rho(X_n, Y_n) < \delta_{in}(\epsilon)}} |\hat{s}_{in}(Y_n) - \hat{s}_{in}(X_n)| D(dX_n, dY_n) \\
& \leq 2A_i \delta_{in}(\epsilon) + 2A_i \eta_i + \int_{\substack{X_n \in \mathcal{E}_{in} \\ Y_n: \rho(X_n, Y_n) < \delta_{in}(\epsilon)}} |\hat{s}_{in}(Y_n) - \hat{s}_{in}(X_n)| D(dX_n, dY_n) \\
& \leq 2A_i \delta_{in}(\epsilon) + 2A_i \eta_i + \epsilon_i \\
& \leq 2A_i \frac{\epsilon}{36A_i^2} + 2A_i \frac{\epsilon}{36A_i^2} + \frac{\epsilon}{18A_i} = \frac{\epsilon}{6A_i}
\end{aligned}$$

For  $n \leq n_0$ , we can pick

$$\delta_{in_0}(\epsilon) = \min\left\{\frac{\epsilon}{36A_i^2}, \min_{n \leq n_0} C_{k_i}(n, \frac{\epsilon}{36A_i^2})\right\}$$

to satisfy

$$\begin{aligned}
d_L(Q_{on}, Q_n) & < \delta_{in}(\epsilon) = \delta_{in_0}(\epsilon) \\
& = \int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) < \frac{\epsilon}{6A_i} ; \forall n \leq n_0
\end{aligned}$$

#### Proof of Lemma 4

To take care of  $n > n_0$ , let  $d_L(Q^0, Q) < \delta_n(\epsilon)$  for  $n > n_0$ , where

$\delta_1 > 0$  is given. Pick again  $D(X_n, Y_n)$  with  $Q_{on}, Q_n$  marginals such that

$$D(X_n, Y_n: \rho(X_n, Y_n) \geq \delta_n(\epsilon_i) \leq \delta_n(\epsilon_i)$$

and write

$$\int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \leq 2A_1 \delta_n(\epsilon_1) +$$

$$+ \int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) ; n > n_0$$

$$X_n, Y_n : \rho(X_n, Y_n) < \delta_n(\epsilon_1)$$

For arbitrary  $Q^0(\cdot)$  there is nothing that can be done at this point. So, let  $Q^0(\cdot)$  be  $m$  dependent. That is, if  $x_1, x_2, \dots, x_n$  is distributed according to  $Q^0(\cdot)$ , then every  $x_i$  depends on only  $m$  preceding and following data.

Then, let  $X_n$  consist of  $X_{ki}$ ;  $i = 1, 2, \dots$  vectors of consecutive data that have  $m$  data gaps between  $X_{ki}$  and  $X_{k,i+1}$ . Then,  $X_{ki}, X_{kj}$ ;  $i \neq j$  are independent by  $Q^0(\cdot)$ . Let  $n = kn_1$  and define the following experimental distribution:

$$n_{X_n}^k(\omega_1, \dots, \omega_k) = \frac{1}{n_1} \{ \# \text{ of } X_k \text{'s with } x_1 < \omega_1, x_2 < \omega_2, \dots, x_k < \omega_k \}$$

The typical sequences  $X_n$  are such that  $n_{X_n}^k(\omega_1, \dots, \omega_k)$  approaches  $Q_k^0(X_k)$  for  $n$  large enough.

Or, given  $\eta_i > 0$ , there is some  $n_0$  and some  $\delta(\eta_i) > 0$  such that

$$\forall n > n_0 \Rightarrow Q^0\{d_{L_\rho}(Q_k^0, n_{X_n}^k) \geq \delta(\eta_i)\} \leq \eta_i$$

Let

$$\mathcal{E}_{ni} = \{X_n : d_{L_\rho}(Q_k^0, n_{X_n}^k) < \delta(\eta_i)\}$$

Then

$$Q^0(\mathcal{E}_{ni}) > 1 - \eta_i ; \forall n > n_0$$

Going back on page 18 we have now

$$\int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \leq 2A_1 \delta_n(\epsilon_1)$$

$$+ \int_{\substack{X_n \in \mathcal{E}_{ni} \\ Y_n: \rho(X_n, Y_n) < \delta_n(\epsilon_1)}} |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) + \int_{\substack{X_n \in \mathcal{E}_{ni}^c \\ Y_n: \rho(X_n, Y_n) < \delta_n(\epsilon_1)}} |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n)$$

$$\leq 2A_1 \delta_n(\epsilon_1) + 2A_1 \eta_1 + \int_{\substack{X_n \in \mathcal{E}_{ni} \\ Y_n: \rho(X_n, Y_n) < \delta_n(\epsilon_1)}} |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n)$$

Let

$$\rho(X_n, Y_n) < \delta_n(\epsilon_1) \Rightarrow d_{L_\rho}(n_{X_n}^k, n_{Y_n}^k) < \zeta_\rho(\delta(\epsilon_1)) ; \forall n > n_0$$

where  $\zeta_\rho$  depends on the measure  $\rho(\cdot, \cdot)$ . Then,

$$\begin{aligned} & \int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \\ & \leq 2A_1 \delta_n(\epsilon_1) + 2A_1 \eta_1 + \int_{\substack{X_n: d_{L_\rho}(Q_k^0, n_{X_n}^k) < \delta(\eta_1) \\ Y_n: d_{L_\rho}(n_{X_n}^k, n_{Y_n}^k) < \zeta_\rho(\delta(\epsilon_1))}} |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \end{aligned}$$

$d_{L_\rho}(\cdot, \cdot)$  satisfying the triangle inequality and, being symmetric, we have:

$$d_{L_\rho}(Q_k^0, n_{Y_n}^k) \leq d_{L_\rho}(Q_k^0, n_{X_n}^k) + d_{L_\rho}(n_{X_n}^k, n_{Y_n}^k)$$

So, from above we have

$$\begin{aligned} & \int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \leq 2A_1 \delta_n(\epsilon_1) + 2A_1 \eta_1 \\ & + \int_{\substack{X_n, Y_n: d_{L_\rho}(Q_k^0, n_{X_n}^k) < \delta(\eta_1) \\ d_{L_\rho}(Q_k^0, n_{Y_n}^k) < \delta(\eta_1) + \zeta_\rho(\delta(\epsilon_1))}} |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \end{aligned}$$

Let  $\hat{s}_{in}(X_n)$  be continuous at  $Q_{ok}$ . That is, given  $\epsilon_1/3 > 0$ , there is some  $\mu(\epsilon_1/3) > 0$  such that

$$\forall X_n : d_{L\rho}(n_{X_n}^k, Q_k^0) < \mu(\epsilon_1/3)$$

$$\Rightarrow |\hat{s}_{in}(X_n) - s_i(Q_k^0)| < \epsilon_1/3$$

where  $s_i(Q_k^0)$  the value that  $\hat{s}_{in}(X_n)$  converges for  $n \rightarrow \infty$ .

Choose:

$$\eta_1 = \frac{\epsilon_1}{6A_1} = \frac{\epsilon}{36A_1^2}; \quad \epsilon_1 = \frac{\epsilon}{6A_1}$$

$$\delta(\eta_1) = \min(\delta(\eta_1 = \frac{\epsilon}{36A_1^2}), \frac{1}{2}\mu(\frac{\epsilon}{18A_1}))$$

$$\delta(\epsilon_1) = \min(\zeta_\rho^{-1}(\frac{1}{2}\mu[\frac{\epsilon}{18A_1}]), \frac{\epsilon}{36A_1^2})$$

Then

$$\begin{aligned} & 2A_1\delta(\epsilon_1) + 2A_1\eta_1 + \int |\hat{s}_{in}(X_n) - \hat{s}_{in}(Y_n)| D(dX_n, dY_n) \\ & X_n, Y_n : d_{L\rho}(Q_k^0, n_{X_n}^k) < \delta(\eta_1) \\ & d_{L\rho}(Q_k^0, n_{Y_n}^k) < \delta(\eta_1) + \zeta_\rho(\delta(\epsilon_1)) \end{aligned}$$

$$\leq \frac{\epsilon}{18A_1} + \frac{\epsilon}{18A_1} + \frac{\epsilon}{18A_1} = \frac{\epsilon}{6A_1}$$

Finally, pick for every  $i$  :

$$\delta(\epsilon) = \min_i \delta(\epsilon_i) =$$

$$= \min_i \min\left\{\frac{\epsilon}{48A_1^2}, \min_{n \leq n_0} c_{k_i}(n, \frac{\epsilon}{48A_1^2}), \zeta_\rho^{-1}(\frac{1}{2}\mu[\frac{\epsilon}{24A_1}])\right\}$$

and then  $d_V(D^0(\hat{S}_n), D(\hat{S}_n)) < \epsilon$ ,

for  $d_{L\rho}(Q^0, Q) < \delta(\epsilon)$ .